

SOME REMARKS ON BAYESIAN PERTURBATION
DIAGNOSTICS AND ROBUSTNESS

by

Seymour Geisser

University of Minnesota

Technical Report #521

November 1988

*This work was sponsored in part by an NIH grant GMS-2527

0. Introduction

To paraphrase the opening remarks of G.E.P. Box (1980) "No [Bayesian] analysis can safely be assumed to be adequate. Perspicacious criticism employing diagnostic checks must therefore be applied."

A Bayesian analysis may depend critically on the modeling assumptions which include prior, likelihood and loss function. While a loss function is presumably a choice made in the context of particular situations, there is no harm and potentially some gain in investigating the effect on an analysis using alternative loss functions. The likelihood is supposed to represent to some approximation the physical process generating the data while the prior reflects subjective views about some of the assumed constructs of this process. Now a likelihood model that has been judged adequate in previous situations similar to a current one is certainly a prime candidate for modeling. However even in such situations the statistician is still obliged to investigate its present adequacy. A way of addressing this problem is to perturb the "standard" model to a greater or lesser degree in potentially conceivable directions to determine the effect of such alterations on the analysis. While for the strict Bayesian the prior is subjective, it is common knowledge how difficult it often is to subject an investigator or even a statistician to an elicitation procedure that convincingly yields an appropriately subjective prior. Hence to perturb an investigator's prior or some standard one that appears appropriate, is also sensible.

1. Types of Perturbation

Even when a standard statistical model has proven adequate in data sets

similar to a current one at hand, one is obliged to consider the effect of perturbing the standard model in one way or another on the analysis especially if graphical procedures indicate the possibility that the standard model may only be marginally adequate.

There are a large number of possible perturbation schemata. A typically useful one is where $w \in \Omega$ an index governing a perturbation schema is a set of hyperparameters. For $X^{(N)} = (X_1, \dots, X_N)$, a set of observables, a rather simple example is

$$f(x^{(N)} | \theta, w) \propto \prod_{j=1}^N \left[1 + \frac{(x_j - \mu)^2}{w \sigma^2} \right]^{-(w+1)/2}, \quad w \geq 1$$

where the standard is $w \rightarrow \infty$, i.e. the normal distribution and the most deviant $w = 1$, the Cauchy distribution.

A second set is exemplified by a mixture e.g.

$$f(x^{(N)} | \theta, w) \propto \prod_{j=1}^N [w f_1(x_j | \alpha) + (1-w) f_2(x_j | \beta)], \quad 0 \leq w \leq 1$$

where say $w = 1$ is the standard and α and β are subsets of θ .

Use of w as an indicator is relevant to situations where w changes the model distribution to varyingly different but known distributional forms not necessarily in the same family. Although this can often be regarded as a special case of either of the first two methods it is best to consider it separately.

A fourth possibility is the use of w as an exclusion indicator i.e. $X^{(N)} = (X_1, \dots, X_N)$ has some standard distributional form under w_0 but for $w \neq w_0$ one or more of the X_i 's have either another distributional form or a completely

unspecifiable distribution. In the former case this could mean for example that an observation's variance differs from the others or more generally that a parameter set not under scrutiny differs for a few of the observations. The latter situation is typically reflected in problems with outliers and aberrant observations that defy satisfactory alternative modeling.

A fifth possibility has to do with what one may term periparametric models. Here $w = w_0$ specifies a standard density while $w \neq w_0$ specifies all model densities $f(x^{(N)}|w)$ that are within a given neighborhood of $f(x^{(N)}|w_0)$ determined by varying w .

A sixth may have to do with possibly inaccurate measurement of the covariates under $w \neq w_0$ or even the actual responses themselves. All of the above have to do essentially with perturbation of the likelihood. Similar remarks may be made regarding the prior $g(\theta|w)$ and combinations of both likelihood and prior. As a typical example the prior could be a mixture e.g.

$$g(\theta|w) = w g_1(\theta) + (1-w)g_2(\theta)$$

with $w = 1$ resulting on the standard $g_1(\theta)$ based on previous information while $g_2(\theta)$ expresses the possibility of another view of the situation. This, for whatever it is worth, results in simpler calculations than having to deal with a likelihood mixture. In particular the use of periparametric perturbation models for additional uncertainty about a "standard" prior seems to be a promising approach especially when combined with a standard likelihood. Here one can examine the extent to which bounds on the "standard" prior can be expanded and still yield moderate sample size robustness, e.g. Lavine (1988).

2. Formal Parametric Analyses

A formal Bayesian framework for a perturbation analysis either for a "relevant" parameter or future observables can be delineated.

For the relevant parameter say θ , we can consider the modeling is such that for a given perturbation index w , the posterior probability function for θ is specified as

$$\mathcal{P}(\theta | x^{(N)}, w) \propto f(x^{(N)} | \theta, w) g(\theta | w)$$

where $g(\theta | w)$ is an assumed prior density for θ conditional on $w \in \Omega$ where $w = w_0$ is the standard.

A loss function

$$L(a, \theta)$$

for taking action $a(X^{(N)}) \in \underline{A}$ upon observing $X^{(N)}$, given θ is the true value, is assumed (the loss function itself may also be perturbed but we shall not consider this possibility in what follows). The average loss

$$\bar{L}_w(a) = \int L(a, \theta) \mathcal{P}(\theta | x^{(N)}, w) d\theta,$$

which depends on w , is now minimized

$$\min_a \bar{L}_w(a) = \bar{L}_w(a_w^*)$$

yielding optimal action a_w^* when w is "true". We then consider the difference in the loss when taking action $a_w^* - a_o^*$, the optimal action under the standard and when $w \neq w_0$ is true. We define the differential loss as

$$d(w) = \bar{L}_w(a_o^*) - \bar{L}_w(a_w^*) \geq 0.$$

One then can examine this loss over a possible range of w to assess its

importance with regard to the action taken under w_0 and in particular

$d^* = \max_{w \in \Omega} d(w)$. We could also assess its local significance by examining $d(w)$ in

a neighborhood about w_0 . In fact if w is a scalar and the second derivative of

$d(w)$ exists and is continuous the calculation of the curvature at $w = w_0$ i.e.

$d''(w_0)$, since $d'(w_0) = 0$, could be rather informative regarding local

perturbations. For example a large curvature would indicate that the actions

taken could be highly sensitive to a slight variation in the standard model.

For a vector w , the matrix of second derivatives will govern the local curvature and one could assess the maximum curvature i.e. in the direction of the normed vector associated with the largest root of the matrix of second derivatives evaluated at the standard $w = w_0$. Cook (1986) has proposed probing local curvature with regard to the displacement of maximized log-likelihoods.

3. Predictive Analysis

We now outline the situation for prediction. The model considered is the joint probability function

$$f(x^{(N)}, x_{(M)}, \theta | w) = f(x_{(M)} | x^{(N)}, \theta, w) f(x^{(N)} | \theta, w) g(\theta | w)$$

whence we obtain

$$f(x_{(M)} | x^{(N)}, w) = \frac{\int f(x^{(N)}, x_{(M)}, \theta | w) d\theta}{\iint f(x^{(N)}, x_{(M)}, \theta | w) d\theta dx_{(M)}}$$

Now assume that $L(a, x_{(M)})$ is the loss incurred in taking action a when observing $x^{(N)}$ given a future realization $x_{(M)}$. The average predictive loss

$$\bar{L}_w(a) = \int L(a, x_{(M)}) f(x_{(M)} | x^{(N)}, w) dx_{(M)}$$

is then minimized

$$\min_a \bar{L}_w(a) = \bar{L}_w(a_w^*)$$

where a_w^* is the optimal action. As before letting $a_{w_0}^* = a^*$, we define the differential loss as $d(w) = \bar{L}_w(a^*) - \bar{L}_w(a_w^*)$ and examine globally $\max_w d(w)$ to determine the possible extent of the maximum effect of the perturbations.

Further in regular cases one can again study locally the maximum curvature which occurs in the direction of the normal vector associated with the largest root of the Hessian matrix, say $d''(w_0)$. If local curvature is appreciable it would appear that the sample is not even robust locally and a review of the standard model is in order. Of course if the perturbed w model is deemed

reasonable one possibility is to define a prior distribution for w and then integrate it out to obtain

$$f(x_{(M)} | x^{(N)}) = \int g(w) f(x_{(M)} | x^{(N)}, w) dw.$$

4. Other Perturbation Diagnostics

Often, we are not in a position to discuss decisions or actions which would necessarily flow from a data set and consequently report either the posterior or predictive distribution itself or some high probability density region for θ or $X_{(M)}$. For reporting the entire posterior distribution the Kullback-Leibler estimative divergence,

$$K(\phi_w, \phi_{w_0}) = E[\log \phi_w - \log \phi_{w_0}],$$

where $\phi_w = \phi(\theta | y^{(N)}, w)$, is a reasonable diagnostic to consider when it exists and is finite, Geisser (1985), Johnson and Geisser (1985), McCulloch (1986) and can be investigated in a variety of paradigms. Similarly for predictive distributions a predictive divergence

$$K(w, w_0) = E[\log f_w - \log f_{w_0}]$$

where $f_w = f(x_{(M)} | x^{(N)}, w)$, will serve as a reasonable diagnostic. Divergences of this sort were introduced by Johnson and Geisser (1982, 1983) for determining influential observations, one of the particular types of perturbation previously mentioned, and were termed predictive influence functions (PIF).

Both, estimative and predictive diagnostics, are most useful in indicating the relative effect of various perturbations.

There may be, however, some difficulty in adequately interpreting globally

$$\max_{w \in \Omega} K(\phi_w, \phi_{w_0}), \text{ or } \max_{w \in \Omega} K(w, w_0)$$

for some of these paradigms.

Another use is to find the direction in which local perturbations have the greatest effect in terms of normal curvature. It can be shown that under suitable regularity conditions that the matrix of second derivatives of $K(\varphi_w, \varphi_{w_0})$ or $K(w, w_0)$ for w a vector of perturbations, say

$$K''_{w=w_0} = I(w_0)$$

where $I(w_0)$ is the Fisher Information matrix for either the posterior or predictive distribution at $w = w_0$, Kullback (1959). The curvature in direction z where $w(t) = w_0 + tz$ and $z'z = 1$ is

$$C_z = z'I(w_0)z$$

so that the maximum curvature C^* is in the direction z^* , the vector associated with the maximum root of $I(w_0)$, where

$$C^* = z^* I(w_0) z^*.$$

An examination of the components of z^* will indicate which ones, namely the larger ones, are those perturbations which relatively most alter the posterior or predictive distribution in terms of the divergence.

Once potentially significant directions are identified, an analysis involving these directions is in order to ascertain whether local departures for them are important enough to vitiate the standard analysis.

The L^1 norm between two densities f and g , favored by Devroye (1987), or the L^2 norm between \sqrt{f} and \sqrt{g} favored by Pitman (1979) as measures of distance between densities can also be used here as diagnostics for the posterior distribution. More generally the Hellinger distance between densities raised to the n -th power

$$H^n = \int |f^{1/n} - g^{1/n}|^n dx$$

yields these as special cases. For the case here with $n=2$ we have

$$H^2(p_w, p_{w_0}) = \int (p_w^{1/2} - p_{w_0}^{1/2})^2 d\theta$$

for posterior densities and for predictive densities,

$$H^2(w, w_0) = \int (f_w^{1/2} - f_{w_0}^{1/2})^2 dx_{(M)},$$

which accords with the norm favored by Pitman. Under suitable smoothness conditions, twice the matrix of second derivatives evaluated at $w=w_0$ when w is a vector, is

$$2H_{w=w_0}^{2''} = I(w_0)$$

again Fisher's Information matrix, which is an indicator of the usefulness of this quantity for local perturbation analysis.

The L^1 norm may also be used. While it is unaffected by any one-to-one transformation as is the divergence and H^2 , it is analytically awkward and does not discriminate between differences of the two densities when the smaller of the two is large or small, as does H^2 and the divergence.

While the divergence and H^2 are as sensible as any measure of how densities differ overall it is difficult to establish a reasonable calibration that different values of the divergence or H^2 entail except in a relative sense. Methods for a more suitably direct interpretation that a statistician, and more to the point an investigator, can readily understand can also be defined but they involve rather specific situations. We now present some of these ways of assessing the robustness in terms of posterior or predictive regions for θ or $X_{(M)}$. One could restrict oneself to perturbations that could matter as determined locally but we shall retain the same notation as before for two reasons. First for convenience in that it is possible that the entire w set may matter and secondly in certain instances one may not be specifically interested

in a local determination. The potential value of the local analysis is the possibility of restricting the dimension of the vector of perturbations to a small set that can more easily be managed by the assessments we now shall propose.

The first method is to assess the robustness of a $1-\alpha$ highest probability density region based on the standard w_0 . Suppose this region denoted by $R_{1-\alpha}(w_0)$ has volume $V(w_0)$ and when perturbed the highest probability density region $R_{1-\alpha}(w)$ has volume $V(w)$. Let $\nu(w)$ be the volume of the intersection of $R_{1-\alpha}(w)$ and $R_{1-\alpha}(w_0)$ as a function of w ,

$$\nu(w) = \text{volume } [R_{1-\alpha}(w) \cap R_{1-\alpha}(w_0)]$$

and let

$$\Gamma_w = \frac{\nu(w)}{M(w)},$$

where $M(w) = \max (V(w), V(w_0))$, be the ratio of the volume of the intersection to whichever is larger the standard or the perturbed for the given w . Then calculate

$$\min_{w \in \Omega} \Gamma_w = \Gamma_w^*$$

which now yields the proportion of the region for the "worst" possible case at a given probability $1-\alpha$. Hence one has an easily interpretable value for assessing the robustness of the data set in terms of a standard analysis involving a $1-\alpha$ region in the presence of presumably anticipated perturbations.

A second method focuses on the use of the standard region's $R_{1-\alpha}(w_0)$ perturbed probability when $w \neq w_0$. Here we use either

$$\Pr[\theta \in R_{1-\alpha}(w_o) | w] = \int_{R_{1-\alpha}(w_o)} p(\theta | y^{(N)}, w) d\theta = 1 - \alpha_w$$

or

$$\Pr[X_{(M)} \in R_{1-\alpha}(w_o) | w] = \int_{R_{1-\alpha}(w_o)} f(x_{(M)} | x^{(N)}, w) dx_{(M)} = 1 - \alpha_w$$

and either

$$\max_{w \in \Omega} |1 - \alpha - (1 - \alpha_w)| = \max_{w \in \Omega} |\alpha_w - \alpha|$$

or

$$\max_{w \in \Omega} \frac{|\alpha_w - \alpha|}{1 - \alpha_w}$$

as easily interpretable values. This second method is most compelling when some specified region is critical to an analysis, e.g. the effect of the perturbation on the calculation of the probability of an observable exceeding some threshold.

In fact as a very simple illustration of this consider $X^{(N)} = (X_1, \dots, X_N)$ a random sample from

$$f(x | \theta, w) = \theta e^{-\theta(x-w)} \quad x \geq w \geq 0$$

and noninformative prior

$$g(\theta) \propto \theta^{-1}.$$

Suppose x_1, \dots, x_d are fully observed realizations and X_{d+1}, \dots, X_N are independently censored at values x_{d+1}, \dots, x_N . We further suppose, as is almost always the case, that

$$x_{(1)} = \min(x_1, \dots, x_d) \leq \min(x_{d+1}, \dots, x_N).$$

The predictive distribution function, Geisser (1982), is then easily calculated to be

$$\Pr[X_{N+1} \leq x | x^{(N)}, w] = 1 - \left(1 + \frac{x - w}{N\bar{x}}\right)^{-d} \quad \begin{matrix} 0 \leq w \leq x_{(1)} \\ w \leq x \end{matrix}$$

Here it is of interest to calculate the probability of a survival threshold

say y

$$\Pr [X_{N+1} > y \mid X^{(N)}, w] = \left(1 + \frac{y - w}{N\bar{x}}\right)^{-d}$$

where the standard say is $w=0$. Of course the divergence and the H^n distances are largely irrelevant for this case but we can easily calculate

$$\max_{0 \leq w \leq x_{(1)}} |\alpha_w(y) - \alpha(y)| = \left(1 + \frac{y - x_{(1)}}{N\bar{x}}\right)^{-d} - \left(1 + \frac{y}{N\bar{x}}\right)^{-d}$$

for a fixed y or conversely for those values of y such that the quantity on the right is no larger than a given value considered negligible with respect to stating a probability for surviving the threshold.

As an example consider the following data reported in Gnedenko et al (1969, p. 176) consisting of a sample of $N=100$ items tested and time to failure recorded for each item until 500 standard time units have elapsed. The recorded failure times for 11 items were: 31, 49, 90, 135, 161, 249, 323, 353, 383, 436, 477. The remaining 89 items survived the test termination time. If interest is focused on the probability of a future item surviving 500 time units then

$$\begin{aligned} \Pr [X_{N+1} > 500 \mid w] &= \left(1 + \frac{500 - w}{47,187}\right)^{-11} \\ &= .891 \quad \text{for } w = 0 \\ &= .897 \quad \text{for } w = x_{(1)} = 31 \end{aligned}$$

Hence

$$\max_{0 \leq w \leq 31} |\alpha_{31}(500) - \alpha(500)| = .006$$

On the other hand one might be interested in that value y such that

$$\Pr[X_{N+1} > y \mid w] = .5$$

Here for $w=0$, $y = 3069$ and for $w = 31$, $y = 3100$ yielding a maximum relative difference of 1%. In passing we also point out here that the maxima for the

divergence and the two norms are $K = \infty$, $H^2 : H^1 = .01$ and are not particularly informative. The divergence indicates only a difference in support while the norms are approximately and exactly twice the probability assigned to the largest interval over which only one of the densities is supported.

More generally, implementation of these methods in other cases could involve the algebraic or numerical calculation of the intersection of two n -dimensional hyperellipsoids which could be quite burdensome for $n > 3$.

Even more complex situations arise where the highest probability density regions are disconnected. Here one may also want to take into consideration the distance from the standard a perturbed and disconnected region is in ordering the diagnostics discussed above, i.e. not only the size of the non-intersecting disconnected region but its distance in some sense from the standard.

5. Acknowledgment

This work was sponsored in part by an NIH grant GMS-25271.

References

- Box, G.E.P.(1981). Sampling and Bayes' inference in scientific modelling and robustness. Journal of the Royal Statistical Society A, 143, 383-430.
- Cook, R.D.(1986). Assessment of local influence (with discussion). Journal of the Royal Statistical Society B, 48, 2, 133-169.
- Devroye, L. (1987). A Course in Density Estimation. Birkhauser.
- Geisser, S. (1982) Aspects of predictive and estimative approaches in the determination of probabilities, Biometrics Supplement: Current Topics in Biostatistics and Epidemiology 38, 1, March, 75-85.
- Geisser, S. (1985) On the predicting of observables: a selective update, in: Bernardo, J.M. et al. (Ed.) Bayesian Statistics 2, (with discussion) 203-230. Amsterdam, North-Holland.
- Gnedenko, B.B., Belyayev, Y.K., and Solov'yev, A.D. (1969). Mathematical Methods of Reliability Theory. New York and London: Academic Press.
- Johnson, W. & Geisser S. (1982) Assessing the predictive influence of observations, in: G. Kallianpur, P.R. Krishnaiah & J.K. Ghosh (Eds) Statistics and Probability Essays in Honor of C.R. Rao, 343-358. Amsterdam, North-Holland.
- Johnson, W. & Geisser, S. (1983) A predictive view of the detection and characterization of influential observations in regression analysis, Journal of American Statistical Association 78, 137-144.
- Johnson, W. & Geisser, S. (1985) Estimative influence measures for the multivariate general linear model, Journal of Statistical Planning and Inference 11, 33-56.
- Kullback, S. (1959). Information theory and Statistics. New York, John Wiley and Sons.
- Lavine, M. (1987). Prior influence in Bayesian Statistics. University of Minnesota Technical Report No. 504
- McCulloch, R. (1986). Local prior influence. University of Minnesota Technical Report No 477.
- Pitman, E.J.G. (1979). Some Basic Theory for Statistical Influence. London, Chapman and Hall.